

FAIR Principles and the Immune Epitope Database

James A. Overton; Knocean Inc.

Randi Vita, Alessandro Sette, Bjoern Peters; La Jolla Institute
Christopher J. Mungall; Lawrence Berkeley National Laboratory


2017-10-25

The Immune Epitope Database

The IEDB by the Numbers

A publicly available database of experiments demonstrating recognition of immune epitopes by adaptive immune receptors.

- ▶ established 2004 by NIAID
- ▶ based at the La Jolla Institute of Allergy and Immunology
- ▶ 18,804 journal articles + direct submissions
- ▶ 1,509,095 assays (individual experiments)
- ▶ up to 400 fields per assay
- ▶ 385,446 epitopes

Epitope	
Epitope ID	717
Chemical Type	Linear peptide
Linear Sequence	ADLMGYIPLVGAPLGGAARA
Starting Position	131
Ending Position	150
Source Molecule Name	polyprotein
Source Accession	BAA03375.1 
Source Organism ID	31647
Source Organism	Hepatitis C virus subtype 1b

Name		John Doe	
Address		123 Main Street, Suite 500, New York, NY 10001	
Phone Number		+1 (212) 555-1234	
Email Address		john.doe@example.com	
Date of Birth		1985-03-15	
Gender		Male	
Marital Status		Single	
Education		Bachelor's Degree	
Employment		Software Engineer	
Languages		English, Spanish	
Interests		Reading, Hiking, Music	
Emergency Contact		Jane Doe (Sister) - +1 (212) 555-5678	
Notes		Client since 2018. Excellent feedback.	

Name		Jane Smith	
Address		456 Elm Street, Apt. 101, Los Angeles, CA 90001	
Phone Number		+1 (310) 555-9876	
Email Address		jane.smith@example.com	
Date of Birth		1990-07-22	
Gender		Female	
Marital Status		Married	
Education		Master's Degree	
Employment		Marketing Specialist	
Languages		English, French	
Interests		Gardening, Travel, Art	
Emergency Contact		Mark Smith (Husband) - +1 (310) 555-4321	
Notes		New client. Initial consultation scheduled for next week.	

Name		Mark Johnson	
Address		789 Oak Avenue, Suite 200, Chicago, IL 60601	
Phone Number		+1 (312) 555-2345	
Email Address		mark.johnson@example.com	
Date of Birth		1978-11-03	
Gender		Male	
Marital Status		Divorced	
Education		High School Graduate	
Employment		Sales Representative	
Languages		English	
Interests		Sports, Fishing, Golf	
Emergency Contact		Sarah Johnson (Daughter) - +1 (312) 555-6789	
Notes		Long-term client. Regular maintenance services.	

Name		Sarah Williams	
Address		321 Pine Street, Suite 100, San Francisco, CA 94102	
Phone Number		+1 (415) 555-3456	
Email Address		sarah.williams@example.com	
Date of Birth		1982-05-18	
Gender		Female	
Marital Status		Single	
Education		Bachelor's Degree	
Employment		Product Manager	
Languages		English, Japanese	
Interests		Cooking, Yoga, Travel	
Emergency Contact		David Williams (Brother) - +1 (415) 555-7890	
Notes		Client since 2020. Very satisfied with services.	

Name		David Brown	
Address		654 Maple Street, Suite 300, Miami, FL 33101	
Phone Number		+1 (305) 555-8901	
Email Address		david.brown@example.com	
Date of Birth		1975-09-07	
Gender		Male	
Marital Status		Married	
Education		Bachelor's Degree	
Employment		Business Development	
Languages		English, Italian	
Interests		Wine Tasting, Soccer, Music	
Emergency Contact		Emily Brown (Wife) - +1 (305) 555-2109	
Notes		Client since 2019. Excellent communication.	

Name		Emily Davis	
Address		987 Cedar Street, Suite 150, Seattle, WA 98101	
Phone Number		+1 (206) 555-4567	
Email Address		emily.davis@example.com	
Date of Birth		1988-12-01	
Gender		Female	
Marital Status		Single	
Education		Master's Degree	
Employment		Data Analyst	
Languages		English, German	
Interests		Reading, Hiking, Photography	
Emergency Contact		Michael Davis (Father) - +1 (206) 555-9012	
Notes		New client. Initial assessment completed.	

Name		Michael Green	
Address		210 Birch Street, Suite 400, Boston, MA 02101	
Phone Number		+1 (617) 555-6543	
Email Address		michael.green@example.com	
Date of Birth		1970-04-25	
Gender		Male	
Marital Status		Married	
Education		Bachelor's Degree	
Employment		Software Engineer	
Languages		English, Russian	
Interests		Chess, Gardening, Travel	
Emergency Contact		Anna Green (Wife) - +1 (617) 555-3210	
Notes		Client since 2017. Regular check-ins.	

Name		Anna White	
Address		543 Elm Street, Suite 250, Denver, CO 80201	
Phone Number		+1 (303) 555-7654	
Email Address		anna.white@example.com	
Date of Birth		1985-08-10	
Gender		Female	
Marital Status		Single	
Education		Bachelor's Degree	
Employment		Marketing Specialist	
Languages		English, Italian	
Interests		Dancing, Art, Travel	
Emergency Contact		Roberto White (Brother) - +1 (303) 555-8901	
Notes		Client since 2021. Excellent feedback.	

Name		Roberto Garcia	
Address		876 Oak Avenue, Suite 100, Phoenix, AZ 85001	
Phone Number		+1 (602) 555-9012	
Email Address		roberto.garcia@example.com	
Date of Birth		1972-06-14	
Gender		Male	
Marital Status		Married	
Education		Bachelor's Degree	
Employment		Sales Representative	
Languages		English, Spanish	
Interests		Fishing, Golf, Travel	
Emergency Contact		Maria Garcia (Wife) - +1 (602) 555-2345	
Notes		Client since 2018. Regular maintenance services.	

Name		Maria Lopez	
Address		109 Pine Street, Suite 300, San Antonio, TX 78201	
Phone Number		+1 (214) 555-3456	
Email Address		maria.lopez@example.com	
Date of Birth		1980-02-28	
Gender		Female	
Marital Status		Single	
Education		Bachelor's Degree	
Employment		Software Engineer	
Languages		English, Spanish	
Interests		Reading, Hiking, Music	
Emergency Contact		Carlos Lopez (Brother) - +1 (214) 555-6789	
Notes		Client since 2020. Excellent communication.	

Name		Carlos Hernandez	
Address		432 Maple Street, Suite 150, Portland, OR 97201	
Phone Number		+1 (503) 555-7890	
Email Address		carlos.hernandez@example.com	
Date of Birth		1978-10-05	
Gender		Male	

FAIR is FAIR

<https://www.force11.org/group/fairgroup/fairprinciples>

To be Findable:

- ▶ F1. (meta)data are assigned a globally unique and persistent identifier
- ▶ F2. data are described with rich metadata (defined by R1 below)
- ▶ F3. metadata clearly and explicitly include the identifier of the data it describes
- ▶ F4. (meta)data are registered or indexed in a searchable resource

To be Accessible:

- ▶ A1. (meta)data are retrievable by their identifier using a standardized communications protocol
 - ▶ A1.1 the protocol is open, free, and universally implementable
 - ▶ A1.2 the protocol allows for an authentication and authorization procedure, where necessary
- ▶ A2. metadata are accessible, even when the data are no longer available

To be Interoperable:

- ▶ I1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.
- ▶ I2. (meta)data use vocabularies that follow FAIR principles
- ▶ I3. (meta)data include qualified references to other (meta)data

To be Reusable:

- ▶ R1. meta(data) are richly described with a plurality of accurate and relevant attributes
 - ▶ R1.1. (meta)data are released with a clear and accessible data usage license
 - ▶ R1.2. (meta)data are associated with detailed provenance
 - ▶ R1.3. (meta)data meet domain-relevant community standards

Summary: FAIR and Square

F1	Persistent IDs
F2	Rich Metadata
F3	Metadata Links to Data
F4	Metadata Indexed
A1	Open Protocols
A2	Metadata Remain Accessible
I1	Open Language
I2	Open Vocabularies
I3	Metadata Links to Metadata
R1	Rich Metadata
R1.1	License
R1.2	Provenance
R1.3	Community

Discussion: A FAIR Shake

Data and Metadata

The distinction between data and metadata is not always clear or helpful.
For the current purposes, the IEDB is a repository of metadata.

Principles and Practises

FAIR principles tell us **what** to do, and sometimes **why** to do it.

They do not tell us **how** to do it.

Shared principles are not enough for interoperability.

FAIR and OBO

The Open Biomedical Ontologies (OBO) community has a set of principles, much in the same spirit as FAIR.

More importantly, OBO also has shared best practises, common infrastructure, open source tooling.

We see an important role for OBO to articulate shared best practises that make the FAIR principles concrete and effective.

My Opinions

	Principle	Practise
F1	Persistent IDs	IRIs and a PURL system
F2	Rich Metadata	RDF(S), OWL, DC, PROV, OBO metadata
A1	Open Protocols	HTTP(S), RDF/OWL in XML, JSON-LD
I1	Language	RDF/OWL
I2	Vocabularies	OBO community
R1.1	License	Creative Commons: CC-BY or CC0
R1.2	Provenance	PROV
R1.3	Community	OBO community

The IEDB: FAIRly Good

F1. Persistent IDs

The IEDB has persistent URLs for assay, reference, epitope:

- ▶ <http://www.iedb.org/assay/1288922>
- ▶ <http://www.iedb.org/reference/1001817>
- ▶ <http://www.iedb.org/epitope/717>

F2. Rich Metadata

Providing rich metadata about experiments is the heart of the IEDB mission.

We can do better!

- ▶ use standardized predicates
- ▶ share richer modelling
- ▶ RDF/OWL representation

F3. Metadata Links to Data

The IEDB links to journal articles by PubMed ID.

We also capture more specific location within the document (as free text): figure X, table Y.

We can do better!

- ▶ standardized, formalized location data
- ▶ link out to PubMed Central structured representations

F4. Metadata Indexed

The IEDB is itself an indexed, searchable resource.

We have also submitted our data to:

- ▶ Biosharing
- ▶ bioCADDIE
- ▶ Wikidata

A1. Open Protocols

The IEDB is:

- ▶ accessible by HTTP, no authentication required
- ▶ published as HTML pages and CSV tables

We can do better!

- ▶ CSV with IRIs and labels
- ▶ RDF/OWL representation as JSON-LD or RDFa



http://8.37.117.76/epitope/123885

NEW TEST



```
epitopes. Includes more than 95% of all
published infectious disease, allergy,
autoimmune, and transplant epitope data.">
17 <meta name="twitter:image"
content="http://www.iedb.org/images/IEDB.png">
18 <meta name="twitter:card" content="summary">
19 <!-- JSON-LD tag -->
20 <script type="application/ld+json">
21 {
22   "http://purl.org/dc/terms/license":
"http://creativecommons.org/licenses/by/4.0/"
23 ,
24 "http://www.w3.org/ns/prov#wasGeneratedBy":
"http://scicrunch.org/browse/resources/SCR_00
6604"
25 }
</script>
```

Unspecified Type

All (1) ▼

Unspecified Type

0 ERRORS 0 WARNINGS



@type

Unspecified Type

http://purl.org/dc/terms/license

<https://creativecommons.org/licenses/by/4.0/>

http://www.w3.org/ns/prov#wasGeneratedBy

https://scicrunch.org/browse/resources/SCR_006604

A2. Metadata Remain Accessible

Published journal articles remain accessible, even if retracted.

The IEDB metadata records will remain accessible in any case.

I1. Knowledge Representation Language

The IEDB uses SQL and publishes in HTML and CSV.

We can do better!

- ▶ CSV exports with IRIs and labels
- ▶ RDF/OWL representation in JSON-LD or RDFa

I2. FAIR Vocabularies: Reuse

The IEDB links to a wide range of resources:

- ▶ UniProt
- ▶ GenBank
- ▶ NCBI Taxonomy
- ▶ PDB
- ▶ IMGT
- ▶ many OBO ontologies
- ▶ more. . .

I2. FAIR Vocabularies: Contributions

The IEDB contributes whenever possible:

- ▶ Ontology for Biomedical Investigations (>300 terms)
- ▶ Disease Ontology (>200 terms)
- ▶ Chemical Entities of Biological Interest (>2400 terms)

I2. FAIR Vocabularies: Development

The IEDB has also developed the MHC Restriction Ontology, following OBO principles:

<http://purl.obolibrary.org/obo/MRO>

I2. FAIR Vocabularies: Public

Sometimes there is no good home for a term that we need: mouse/rat strains, post-translational modifications of proteins, etc.

We put these terms in an application ontology ONTIE, then migrate to better homes when available.

ONTIE is now public with PURLs and RDF/OWL representations in multiple formats:

https://ontology.iedb.org/ontology/ONTIE_0002032

ONTIE:0002032 Polymerase basic protein 2 (Influenza A virus)

https://ontology.iedb.org/ontology/ONTIE_0002032

- **type:** owl:Class
- **label:** Polymerase basic protein 2 (Influenza A virus)
- **alternative term:** PB2
- **alternative term:** RNA-directed RNA polymerase subunit P3
- **IEDB alternative term:** Polymerase basic protein 2
- **template:** protein class
- **ONTIE domain:** protein
- **protein label:** Polymerase basic protein 2
- **protein taxon:** Influenza A virus

Other formats: [Turtle \(ttl\)](#), [JSON-LD \(json\)](#), [TSV \(tsv\)](#).

I2. FAIR Vocabularies: Improvements

We can do better!

We want to do a better job of explaining to users how vocabularies are used.

One example: our assays often compose a method with a GO process.

13. Metadata Link to Metadata

The IEDB links to many resources (12), which link to other resources, which link to other resources, . . .

R1.1. License

The IEDB was failing this principle.

The IEDB's content is available under the Creative Commons Attribution 4.0 International license. We are making this clear with human- and machine-readable annotations.

R1.2. Provenance

The IEDB uses PubMed IDs to link to the original article.

We are using PROV and adding metadata about our generation of the metadata from an original source.

R1.3. Community

We are careful to use standards that are appropriate for the immunology community.

Summary: Playing FAIR

- ▶ F3: Standardize identification of journal parts (Figures / Tables)
- ▶ F4: Add IEDB metadata to Biosharing, Biocaddie and Wikidata
- ▶ A1: Provide machine actionable representation of IEDB assay level data
- ▶ I1: Represent the IEDB data in RDF/OWL
- ▶ I2: Make all links to external vocabularies explicit
- ▶ I2: Make all internal vocabularies public via ONTIE and link to them
- ▶ R1.1: Include licensing information with the IEDB records
- ▶ R1.2: Include provenance information regarding IEDB curation

More than FAIR

More than FAIR

FAIR is not enough.

Interoperability requires shared best practises, compatible implementations.

Full interoperability also requires shared modelling patterns – we must represent similar things in similar ways.

We see a role for OBO here, and have started coordinated modelling of immunology datasets from the IEDB, BRCs, and ImmPort.

Acknowledgements

Funding is provided by The National Institutes of Health

HHSN272201200010C

FAIR Thee Well